

*On the Theory of Correlation for any Number of Variables,  
treated by a New System of Notation.*

By G. UDNY YULE, Newmarch Lecturer on Statistics, University College,  
London.

(Communicated by Professor O. Henrici, F.R.S. Received January 25,—  
Read February 28, 1907.)

1. The systems of notation hitherto used by writers on the theory of correlation are somewhat unsatisfactory when many variables are involved. In the present paper a new notation is proposed which is simple, definite, and quite general, thus very greatly facilitating the treatment of the subject. The majority of the results given in the sequel were, in fact, first suggested by the notation itself.

2. Let  $x_1 x_2 \dots x_n$  denote deviations in the values of the  $n$  variables from their respective arithmetic means. Then the regression equation may be written :—

$$x_1 = b_{12 \cdot 34 \dots n} x_2 + b_{13 \cdot 24 \dots n} x_3 + \dots + b_{1n \cdot 23 \dots n-1} x_n. \quad (1)$$

In this notation the suffix of each regression coefficient completely defines it. The first subscript gives the dependent variable, the second the variable of which the given regression is the coefficient, and the subscripts after the period show the remaining independent variables which enter into the equation. It is convenient to distinguish the subscripts before and after the period as "primary" and "secondary" subscripts respectively. The order in which the secondary subscripts are arranged is indifferent, but the order of the two primary subscripts is material; e.g.,  $b_{12 \cdot 3 \dots n}$  and  $b_{21 \cdot 3 \dots n}$  denote two quite distinct coefficients. A coefficient with  $p$  secondary subscripts may be termed a regression of the  $p$ th order, the total regressions  $b_{12}$ ,  $b_{13}$ ,  $b_{23}$ , etc., being thus regarded as of order zero.

3. The correlation-coefficients may be distinguished by subscripts in precisely the same manner. Thus the correlation  $r_{12 \cdot 34 \dots n}$  is defined by the relation

$$r_{12 \cdot 34 \dots n} = (b_{12 \cdot 34 \dots n} \cdot b_{21 \cdot 34 \dots n})^{\frac{1}{2}}. \quad (2)$$

In the case of the correlations, the order of both primary and secondary subscripts is indifferent. A correlation with  $p$  secondary subscripts may be termed a correlation of order  $p$ , the total correlations  $r_{12}$ ,  $r_{13}$ ,  $r_{23}$ , etc., being regarded as of order zero.

4. If the regressions in equation (1) be determined as usual by the method

of least squares, the difference between  $x_1$  and the expression on the right, for any observed set of values of  $x_1 x_2 \dots x_n$ , may be denoted by  $x_{1 \cdot 23 \dots n}$ : that is

$$x_{1 \cdot 23 \dots n} = x_1 - b_{12 \cdot 34 \dots n} x_2 - b_{13 \cdot 24 \dots n} x_3 - \dots - b_{1n \cdot 23 \dots n-1} x_n. \quad (3)$$

Such a residual, or deviation, denoted by a symbol with  $p$  secondary subscripts may be termed a deviation of the  $p$ th order,  $x_1 x_2 \dots x_n$  being regarded as deviations of order zero.

5. Finally, the standard deviation  $\sigma_{1 \cdot 23 \dots n}$  is defined as given by the relation

$$N \cdot \sigma^2_{1 \cdot 23 \dots n} = \sum (x_{1 \cdot 23 \dots n}^2), \quad (4)$$

$N$  being the number of observations. If the standard deviation be denoted by a symbol with  $p$  secondary subscripts, it is of the  $p$ th order, the total standard deviations being regarded as of order zero.

6. In terms of this notation, the normal equations from which the regressions are determined may be very briefly written, in the form

$$\begin{aligned} \sum (x_2 \cdot x_{1 \cdot 23 \dots n}) &= \sum (x_3 \cdot x_{1 \cdot 23 \dots n}) = \dots \\ &= \sum (x_n \cdot x_{1 \cdot 23 \dots n}) = 0. \end{aligned} \quad (5)$$

That is to say, we have the general theorem: "The product-sum of any deviation of order zero with any deviation of higher order is zero, provided the subscript of the former occur amongst the secondary subscripts of the latter."

7. It follows that the product-sum of any two deviations of the same order, with the same secondary suffixes, is unaltered by omitting any or all of the secondary subscripts of either and, conversely, the product-sum of any deviation of order  $p$  with a deviation of order  $p+q$ , the  $p$  subscripts being the same in each case, is unaltered by adding to the secondary subscripts of the former any or all of the  $q$  additional subscripts of the latter, for we have by § 6:—

$$\begin{aligned} \sum (x_{1 \cdot 34 \dots n} x_{2 \cdot 34 \dots n}) &= \sum (x_{1 \cdot 34 \dots n}) (x_2 - b_{23 \cdot 4 \dots n} x_3 - \dots - b_{2n-3 \dots n-1} x_n) \\ &= \sum (x_{1 \cdot 34 \dots n} x_2). \end{aligned}$$

Similarly

$$\sum (x_{1 \cdot 34 \dots n} x_{2 \cdot 34 \dots n-1}) = \sum (x_{1 \cdot 34 \dots n} x_2),$$

and so on. Therefore, quite generally,

$$\begin{aligned} \sum (x_{1 \cdot 34 \dots n} x_{2 \cdot 34 \dots n}) &= \sum (x_{1 \cdot 34 \dots n} x_{2 \cdot 34 \dots n-1}) = \dots \\ &\dots = \sum (x_{1 \cdot 34 \dots n} x_2). \end{aligned} \quad (6)$$

8. It follows from § 7 as a corollary from § 6 that the product-sum of any two deviations is zero if all the subscripts of the one are contained among the secondary subscripts of the other.

These theorems (§§ 6—8) give the key to simple deductions of many results in the theory of multiple correlation.

9. We have from the last section and § 7,

$$\begin{aligned} 0 &= \Sigma(x_{2.34...n}x_{1.234...n}) \\ &= \Sigma(x_{2.34...n})(x_1 - b_{12.34...n}x_2 - \text{terms in } x_3 \text{ to } x_n) \\ &= \Sigma(x_1x_{2.34...n}) - b_{12.34...n}\Sigma(x_2x_{2.34...n}) \\ &= \Sigma(x_{1.34...n}x_{2.34...n}) - b_{12.34...n}\Sigma(x_{2.34...n}^2). \end{aligned}$$

That is

$$b_{12.34...n} = \frac{\Sigma(x_{1.34...n}x_{2.34...n})}{\Sigma(x_{2.34...n}^2)}. \quad (7)$$

But this is the value that would have been obtained by taking a regression equation of the form

$$x_{1.34...n} = b_{12.34...n}x_{2.34...n},$$

and determining  $b_{12.34...n}$  by the method of least squares. That is to say,  $b_{12.34...n}$  may be regarded, quite generally and without any reference to the form of the frequency distribution, as the regression of  $x_{1.34...n}$  on  $x_{2.34...n}$ . It follows at once from the definition (3) that  $r_{12.34...n}$  may be regarded as the correlation between  $x_{1.34...n}$  and  $x_{2.34...n}$ , and from (4) that we may write

$$b_{12.34...n} = r_{12.34...n} \frac{\sigma_{1.34...n}}{\sigma_{2.34...n}}. \quad (8)$$

All the relations, in fact, that hold good between deviation-sums, standard deviations, regressions and correlations of order zero, are also valid between deviation-sums, standard deviations, regressions and correlations of any high order.

10. This result is of some importance as regards the interpretation of partial correlations and regressions. In the case of normal correlation there is no difficulty in assigning a meaning to these constants, as the regression is strictly linear, and the partial correlations and regressions are the same for all types of the variables. But in the general case this is not so, and although I showed, in a previous discussion of the question,\* that the values assigned to the partial regressions on the assumption of normal correlation are the "least square" values and, consequently, that the partial correlation retains an "average significance," I could not prove that it remains an actual correlation between determinate variables. The above theorem completes the work in this respect. If, with three variables  $x_1$ ,  $x_2$ , and  $x_3$ , for example, the two regressions  $b_{13}$  and  $b_{23}$  be determined in the ordinary way, and then the residuals  $x_{1.3} = x_1 - b_{13}x_3$ ,  $x_{2.3} = x_2 - b_{23}x_3$  be calculated for all sets of observations  $x'_1 x'_2 x'_3$ ,  $x''_1 x''_2 x''_3$ , etc., the correlation between  $x_{1.3}$  and  $x_{2.3}$  is  $r_{12.3}$ . A similar interpretation holds for any greater number of variables.

\* 'Roy. Soc. Proc.' vol. 60 (1897), p. 477; 'Roy. Stat. Soc. Journ.', vol. 60 (1897), p. 812.

Such a relation would not, of course, afford a practical method of calculating the partial coefficients, as the arithmetic would be extremely lengthy.

11. Any standard deviation of order  $p$  may be expressed in terms of a standard deviation of order  $p-1$  and a correlation of order  $p-1$ . For we have, using the theorems of §§ 6 and 7,

$$\begin{aligned}\Sigma(x_{1 \cdot 23 \dots n}^2) &= \Sigma(x_{1 \cdot 23 \dots n-1} x_{1 \cdot 23 \dots n}) \\ &= \Sigma(x_{1 \cdot 23 \dots n-1})(x_1 - b_{1n \cdot 23 \dots n-1} x_n - \text{terms in } x_2 \text{ to } x_{n-1}) \\ &= \Sigma(x_{1 \cdot 23 \dots n-1}^2) - b_{1n \cdot 23 \dots n-1} \Sigma(x_{1 \cdot 23 \dots n-1} x_{n \cdot 23 \dots n-1});\end{aligned}$$

or, dividing through by the number of observations,

$$\begin{aligned}\sigma_{1 \cdot 23 \dots n}^2 &= \sigma_{1 \cdot 23 \dots n-1}^2 (1 - b_{1n \cdot 23 \dots n-1} r_{1n \cdot 23 \dots n-1}) \\ &= \sigma_{1 \cdot 23 \dots n-1}^2 (1 - r_{1n \cdot 23 \dots n-1}^2).\end{aligned}\quad (9)$$

The form of this relation is the same as that of the familiar relation between a standard deviation of the first order and a standard deviation of order zero, with the secondary subscripts  $23 \dots n-1$  added throughout. It is clear from (9) that  $r_{1n \cdot 23 \dots n-1}$  cannot be numerically greater than unity. It also follows at once that if we have been estimating  $x_1$  from  $x_2, x_3 \dots x_{n-1}, x_n$  will not increase the accuracy of estimate unless  $r_{1n \cdot 23 \dots n-1}$  (not  $r_{1n}$ ) differ from zero.\*

12. In equation (9) the subscript  $n$  is eliminated from the suffix of  $\sigma_{1 \cdot 23 \dots n}$ , and it is obvious that any other subscript can be eliminated in the same way. Therefore we must have

$$\begin{aligned}\sigma_{1 \cdot 34 \dots n}^2 (1 - r_{12 \cdot 34 \dots n}^2) &= \sigma_{1 \cdot 24 \dots n}^2 (1 - r_{13 \cdot 24 \dots n}^2) \\ &= \dots = \sigma_{1 \cdot 23 \dots n-1}^2 (1 - r_{1n \cdot 23 \dots n-1}^2).\end{aligned}\quad (10)$$

Further, we have

$$\begin{aligned}\sigma_{1 \cdot 23 \dots n-1}^2 &= \sigma_{1 \cdot 23 \dots n-2}^2 (1 - r_{1(n-1) \cdot 23 \dots n-2}^2) \\ \sigma_{1 \cdot 23 \dots n-2}^2 &= \sigma_{1 \cdot 23 \dots n-3}^2 (1 - r_{1(n-2) \cdot 23 \dots n-3}^2),\end{aligned}$$

and so on; so that

$$\sigma_{1 \cdot 23 \dots n}^2 = \sigma_1^2 (1 - r_{12}^2) (1 - r_{13 \cdot 2}^2) (1 - r_{14 \cdot 23}^2) \dots (1 - r_{1n \cdot 23 \dots n-1}^2). \quad (11)$$

This is an extremely convenient expression for arithmetical use, as illustrated later. A complete check on the arithmetic is obtained by eliminating the secondary subscripts in a different, say the inverse, order, *i.e.*, by using the result—

$$\sigma_{1 \cdot 23 \dots n}^2 = \sigma_1^2 (1 - r_{1n}^2) (1 - r_{1(n-1) \cdot n}^2) (1 - r_{1(n-2) \cdot n(n-1)}^2) \dots (1 - r_{12 \cdot 34 \dots n}^2). \quad (12)$$

\* Cf. proofs for cases of 3 and 4 variables previously given (*loc. cit.* in previous note).

13. Any regression of order  $p$  may be expressed in terms of regressions of order  $p - 1$ . For we have

$$\begin{aligned}\Sigma(x_{1.34\dots n}x_{2.34\dots n}) &= \Sigma(x_{1.34\dots n-1}x_{2.34\dots n}) \\ &= \Sigma(x_{1.34\dots n-1})(x_2 - b_{2n.34\dots n-1}x_n - \text{terms in } x_3 \text{ to } x_{n-1}) \\ &= \Sigma(x_{1.34\dots n-1}x_{2.34\dots n-1}) - b_{2n.34\dots n-1}\Sigma(x_{1.34\dots n-1}x_{n.34\dots n-1}).\end{aligned}$$

That is, replacing  $b_{2n.34\dots n-1}$  by  $b_{n2.34\dots n-1} \times \sigma^2_{2.34\dots n-1}/\sigma^2_{n.34\dots n-1}$ ,

$$b_{12.34\dots n}\sigma^2_{2.34\dots n} = b_{12.34\dots n-1}\sigma^2_{2.34\dots n-1} - b_{1n.34\dots n-1}b_{n2.34\dots n-1}\sigma^2_{2.34\dots n-1}.$$

Therefore, by equation (9),

$$b_{12.34\dots n} = \frac{b_{12.34\dots n-1} - b_{1n.34\dots n-1} \cdot b_{n2.34\dots n-1}}{1 - b_{1n.34\dots n-1}b_{n1.34\dots n-1}}. \quad (13)$$

But this is simply the expression for  $b_{12.n}$  in terms of  $b_{12}$ ,  $b_{1n}$ ,  $b_{n1}$ , and  $b_{n2}$ , with the subscripts  $34\dots n-1$  added throughout. Therefore  $b_{12.34\dots n}$  may be regarded as the partial regression of  $x_{1.34\dots n-1}$  on  $x_{2.34\dots n-1}$ ,  $x_{n.34\dots n-1}$  being given. As any other secondary subscript might have been eliminated in lieu of  $n$ , we can also regard it as the partial regression of  $x_{1.45\dots n}$ , on  $x_{2.45\dots n}$ ,  $x_{3.45\dots n}$  being given, and so on.

14. Equation (13) may be written in terms of the correlations:—

$$b_{12.34\dots n} = \frac{r_{12.34\dots n-1} - r_{1n.34\dots n-1}r_{2n.34\dots n-1}}{1 - r_{1n.34\dots n-1}^2} \frac{\sigma_{1.34\dots n-1}}{\sigma_{2.34\dots n-1}}.$$

Hence, writing down the similar expression for  $b_{21.34\dots n}$ , and taking the square root of the product,

$$r_{12.34\dots n} = \frac{r_{12.34\dots n-1} - r_{1n.34\dots n-1}r_{2n.34\dots n-1}}{(1 - r_{1n.34\dots n-1}^2)^{\frac{1}{2}}(1 - r_{2n.34\dots n-1}^2)^{\frac{1}{2}}}. \quad (14)$$

This is, similarly, the expression for  $r_{12.n}$  in terms of  $r_{12}$ ,  $r_{1n}$ , and  $r_{2n}$ , with the secondary subscripts  $34\dots n-1$  added throughout, and accordingly  $r_{12.34\dots n}$  may be regarded as the partial correlation between  $x_{1.34\dots n-1}$  and  $x_{2.34\dots n-1}$ ,  $x_{n.34\dots n-1}$  being given, and so on, as for the regression.

15. It is clear that equations (13) and (14) imply a series of relations between correlations or regressions of orders less than  $n - 2$  with  $n$  variables, for all the expressions obtained by eliminating  $34\dots n$  in turn from the secondary subscripts of the constant on the left must be equal to each other. Further, every coefficient of the  $p$ th order can be expressed in terms of the coefficients of the  $(p - 1)$ th order in  $p$  different ways, by eliminating each of the  $p$  secondary subscripts in turn. This enables an absolute check to be kept on the arithmetic by calculating each coefficient in at least two distinct ways.

16. By the use of equation (14), the work of calculating correlation coefficients of higher orders is rendered quite simple and straightforward. The use of equation (13) for calculating the regressions is comparatively

clumsy, however: when the correlations have been found, it is best to work out the standard deviations by equation (11), and then the regressions are given at once by (8). The following data, taken from a discussion of pauperism,\* will serve as an arithmetical illustration, the variables being the percentage changes during a decade in the poor-law unions of England in: (1) the percentage of the population in receipt of poor-law relief; (2) the ratio of the numbers given relief out-doors to one indoors (in the workhouse); (3) the proportion of aged (over 60) in the population; (4) the population itself. The values of the correlations of order zero are given in Table I, and the logarithms of  $(1 - r^2)^{\frac{1}{2}}$ , required in the calculations, are entered in the third column. These coefficients are next grouped in sets of three, one set to each possible group of three variables, as in the second column of Table II, and the coefficients of the first order are then calculated from (14). For convenience in calculating the coefficients of the second order, the values of  $\log(1 - r^2)^{\frac{1}{2}}$  are again entered in the last column.

Table I.

Correlation coefficient.		$\log \sqrt{1 - r^2}$ .
12	+ 0.52	1.93154
13	+ 0.41	1.96003
14	- 0.14	1.99570
23	+ 0.49	1.94038
24	+ 0.23	1.98820
34	+ 0.25	1.98598

Table II.

Correlation coefficient (zero order).	Product term of numerator.	Numerator.	Correlation coefficient (first order).	$\log \sqrt{1 - r^2}$ .
12	+ 0.52	+ 0.2009	+ 0.3191	12.3
13	+ 0.41	+ 0.2548	+ 0.1552	13.2
23	+ 0.49	+ 0.2132	+ 0.2768	23.1
12	+ 0.52	- 0.0322	+ 0.5522	12.4
14	- 0.14	+ 0.1196	- 0.2596	14.2
24	+ 0.23	- 0.0728	+ 0.3028	24.1
13	+ 0.41	- 0.0350	+ 0.4450	13.4
14	- 0.14	+ 0.1025	- 0.2425	14.3
34	+ 0.25	- 0.0574	+ 0.3074	34.1
23	+ 0.49	+ 0.0575	+ 0.4325	23.4
24	+ 0.23	+ 0.1225	+ 0.1075	24.3
34	+ 0.25	+ 0.1127	+ 0.1373	34.2

\* 'Roy. Stat. Soc. Journ.', vol. 62 (1899), p. 249.

Table III.

Correlation coefficient (first order).		Product term of numerator.	Numerator.	Correlation coefficient (second order).		$\log \sqrt{1 - r^2}$ .
12·4	+0·5731	+0·2131	+0·3600	12·34	+0·458	1·89774
13·4	+0·4642	+0·2630	+0·2012	13·24	+0·276	1·96559
23·4	+0·4590	+0·2660	+0·1930	23·14	+0·266	1·96814
12·3	+0·4014	-0·0350	+0·4364	12·34	+0·458	—
14·3	-0·2746	+0·0511	-0·3257	14·23	-0·359	1·94007
24·3	+0·1274	-9·1102	+0·2376	24·13	+0·270	1·96713
13·2	+0·2084	-0·0505	+0·2589	13·24	+0·276	—
14·2	-0·3123	+0·0837	-0·3460	14·23	-0·359	—
34·2	+0·1618	-0·0651	+0·2269	34·12	+0·244	1·97333
23·1	+0·3553	+0·1219	+0·2334	23·14	+0·266	—
24·1	+0·3580	+0·1209	+0·2371	24·13	+0·270	—
34·1	+0·3404	+0·1272	+0·2132	34·12	+0·244	—

The first order coefficients, from Table II, are then regrouped according to the same primary subscripts as in Table I, and the work repeated precisely as before, as in Table III, but each coefficient of the second order is automatically calculated by this process in two ways and the work thus checked. Small errors introduced by the non-retention of insignificant figures may, of course, prevent complete agreement to the last place of decimals, and for this reason the coefficients of the first order were evaluated to four figures, although only three were required for the final result. In order to obtain the regression equation between changes in pauperism and changes in the three remaining variables, we require the three regressions  $b_{12\cdot34}$ ,  $b_{13\cdot24}$ , and  $b_{14\cdot23}$  and, accordingly, must obtain the six standard deviations,  $\sigma_{1\cdot34}$ ,  $\sigma_{2\cdot34}$ ,  $\sigma_{1\cdot24}$ ,  $\sigma_{3\cdot24}$ ,  $\sigma_{1\cdot23}$ ,  $\sigma_{4\cdot23}$ .

These are readily calculated and checked by means of the equations of the form—

$$\begin{aligned}\sigma_{1\cdot34} &= \sigma_1 (1 - r^2_{13})^{\frac{1}{2}} (1 - r^2_{14\cdot3})^{\frac{1}{2}} \\ &= \sigma_1 (1 - r^2_{14})^{\frac{1}{2}} (1 - r^2_{13\cdot4})^{\frac{1}{2}}\end{aligned}$$

given  $\sigma_1 = 29\cdot2$ ,  $\sigma_2 = 41\cdot7$ ,  $\sigma_3 = 5\cdot5$ ,  $\sigma_4 = 23\cdot8$ ; and the values found are:—

$$\begin{aligned}\sigma_{1\cdot34} &= 25\cdot61, & \sigma_{1\cdot24} &= 27\cdot63, & \sigma_{1\cdot23} &= 24\cdot39, \\ \sigma_{2\cdot34} &= 36\cdot06, & \sigma_{3\cdot24} &= 4\cdot73, & \sigma_{4\cdot23} &= 22\cdot86.\end{aligned}$$

Hence, from the equations of the form

$$b_{12\cdot34} = r_{12\cdot34} \frac{\sigma_{1\cdot34}}{\sigma_{2\cdot34}},$$

we have

$$b_{12\cdot34} = +0\cdot325, \quad b_{13\cdot24} = +1\cdot383, \quad b_{14\cdot23} = -0\cdot383.$$

That is, the regression equation between changes in pauperism and changes in the other factors considered is

$$x_1 = 0.325x_2 + 1.383x_3 - 0.383x_4.$$

To complete the work, we may calculate  $\sigma_{1.234}$ , the standard error made in estimating  $x_1$  from  $x_2$ ,  $x_3$ , and  $x_4$  by the above equation. The value is

$$\begin{aligned}\sigma_{1.234} &= \sigma_1 (1 - r_{12}^2)^{\frac{1}{2}} (1 - r_{13.2}^2)^{\frac{1}{2}} (1 - r_{14.23}^2)^{\frac{1}{2}} \\ &= \sigma_1 (1 - r_{14}^2)^{\frac{1}{2}} (1 - r_{13.4}^2)^{\frac{1}{2}} (1 - r_{12.34}^2)^{\frac{1}{2}} \\ &= 22.8.\end{aligned}$$

17. If, in accordance with the notation used for elementary cases in the paper already referred to, and that in a recent note by Mr. R. H. Hooker and myself,\* we write

$$\sigma^2_{1.23...n} = \sigma^2_1 (1 - R^2_{1(23...n)}), \quad (15)$$

$R_{1(23...n)}$  may be regarded as a coefficient of correlation between  $x_1$  and the expression

$$e_{1.23...n} = b_{12.34...n}x_2 + b_{13.24...n}x_3 + \dots + b_{1n.23...n-1}x_n. \quad (16)$$

The value of  $R$  is accordingly a useful datum, as indicating how closely  $x_1$  can be expressed in terms of a linear function of  $x_2x_3 \dots x_n$ . It may be readily calculated either direct from the equation

$$1 - R^2_{1(23...n)} = (1 - r_{12}^2)(1 - r_{13.2}^2) \dots (1 - r_{1n.23...n-1}^2), \quad (17)$$

or from the value of  $\sigma_{1.23...n}$  and  $\sigma_1$ , if previously obtained.

It is obvious from (17) that, since every bracket on the right is not greater than unity,

$$1 - R^2_{1(23...n)} \leq 1 - r_{12}^2.$$

Hence  $R_{1(23...n)}$  cannot be numerically less than  $r_{12}$ . For the same reason, rewriting (17) in every possible form,  $R_{1(23...n)}$  cannot be numerically less than  $r_{12}$ ,  $r_{13}, \dots, r_{1n}$ , i.e., any one of the possible constituent coefficients of order zero. Further, for similar reasons,  $R_{1(23...n)}$  cannot be numerically less than any possible constituent coefficient of any higher order. That is to say,  $R_{1(23...n)}$  is not less than the greatest of all the possible constituent coefficients of all orders, and is usually, though not always, markedly greater. Thus in the illustration of § 16, the value of  $R_{1(234)}$  is 0.626, and the greatest correlation coefficient is  $r_{12.34} = 0.458$ . The sign of  $R$  is necessarily positive, for a positive increment in  $x_1$  obviously corresponds on the average to a positive increment in  $e_{1.23...n}$ . More definitely, the standard deviation of  $e_{1.23...n}$  is  $\sigma_1 R_{1(23...n)}$ , and the regression of  $x_1$  on  $e_{1.23...n}$  is therefore + 1.

Seeing that  $\sigma^2_{1.23...n} = \sigma^2_1 (1 - R^2_{1(23...n)})$ , and that  $\sigma_{1.23...n}$  is a minimum, we may, alternatively, regard the values of the regressions as determined by the

\* 'Roy. Stat. Soc. Journal,' vol. 59 (1906), p. 197.

condition that the correlation between  $x_1$  and  $x_{1.23\dots n}$ , viz.,  $R_{1(23\dots n)}$ , shall be a maximum.

18. It is obvious that equations (13) and (14) imply relations of an inverse kind, expressing coefficients of a lower order in terms of those of a higher order. Using the same method of expansion as in previous cases, we have

$$\begin{aligned} 0 &= \Sigma (x_{1.23\dots n} x_{2.34\dots n-1}) \\ &= \Sigma (x_1 x_{2.34\dots n-1}) - b_{12.34\dots n} \Sigma (x_2 x_{2.34\dots n-1}) \\ &\quad - b_{1n.23\dots n-1} \Sigma (x_n x_{2.34\dots n-1}). \end{aligned}$$

That is

$$b_{12.34\dots n-1} = b_{12.34\dots n} + b_{1n.23\dots n-1} b_{n2.34\dots n-1}.$$

But by interchanging the suffixes, viz., 1 for  $n$  and  $n$  for 1,

$$b_{n2.34\dots n-1} = b_{n2.13\dots n-1} + b_{n1.23\dots n-1} b_{12.34\dots n-1}.$$

Substituting this value of  $b_{n2.34\dots n-1}$  in the first equation and simplifying,

$$b_{12.34\dots n-1} = \frac{b_{12.34\dots n} + b_{1n.23\dots n-1} b_{n2.13\dots n-1}}{1 - b_{1n.23\dots n-1} b_{n1.23\dots n-1}}. \quad (18)$$

This is the required equation for the regressions. The similar equation for the correlations is obtained at once by writing down the corresponding expression for  $b_{21.34\dots n-1}$  and taking the square root

$$r_{12.34\dots n-1} = \frac{r_{12.34\dots n} + r_{1n.23\dots n-1} r_{2n.13\dots n-1}}{(1 - r_{1n.23\dots n-1}^2)^{\frac{1}{2}} (1 - r_{2n.13\dots n-1}^2)^{\frac{1}{2}}}. \quad (19)$$

19. The general principle that any equation subsisting between such statistical constants as correlations, regressions, and standard deviations holds good for all secondary subscripts, applies also to the equation (3), which expresses the individual deviation of order  $p$  in terms of deviations of order zero. That is to say, we have, quite generally,  $k$  being any subscript or collection of subscripts,

$$x_{1.2\dots kn} = x_{1.k} - b_{12.3\dots kn} x_{2.k} - \dots - b_{1n.2\dots k(n-1)} x_{n.k}. \quad (20)$$

For, if  $l$  be any one of the subscripts included under  $k$ , and if  $m$  denote the remaining subscripts, on expanding both sides of the equation in terms of deviations of order zero, the coefficients of  $x_1, x_2, \dots x_n$  are the same. The coefficients of  $x_l$  are equal if

$$b_{1l.2\dots mn} = b_{1l.m} - b_{2l.m} \cdot b_{12.3\dots kn} - \dots - b_{nl.m} \cdot b_{1n.2\dots k(n-1)}.$$

But, replacing the regressions  $b_{1l.m}, b_{2l.m}, \dots b_{nl.m}$  by product sums, this reduces to

$$\Sigma (x_{l.m} \cdot x_{1.2\dots lm}) = 0,$$

which is true by § 8, whether  $m$  denote a single subscript or an aggregate, or is absent, and equation (20) is accordingly correct. Remembering that

$$b_{12 \cdot 3 \dots kn} = r_{12 \cdot 3 \dots kn} \frac{\sigma_{1 \cdot 3 \dots kn}}{\sigma_{2 \cdot 3 \dots kn}} = r_{12 \cdot 3 \dots kn} \frac{\sigma_{1 \cdot 23 \dots kn}}{\sigma_{2 \cdot 13 \dots kn}},$$

the equation may also be written in the useful form—

$$\frac{x_{1 \cdot 2 \dots kn}}{\sigma_{1 \cdot 2 \dots kn}} = \frac{x_{1 \cdot k}}{\sigma_{1 \cdot 2 \dots kn}} - r_{12 \cdot 3 \dots kn} \frac{x_{2 \cdot k}}{\sigma_{2 \cdot 13 \dots kn}} - \dots - r_{1n \cdot 2 \dots k(n-1)} \frac{x_{n \cdot k}}{\sigma_{n \cdot 12 \dots k(n-1)}}. \quad (21)$$

20. In all the preceding sections no assumption of any kind has been made with respect to the form of the distribution of frequency, but the results may, of course, be applied to the special case of the normal distribution. Let  $y_{12 \dots n}$  denote the value of the normal function for the combination of deviations  $x_1, x_2, \dots, x_n$ , and  $y'_{12 \dots n}$  the value of the function when all deviations are zero, then we may write

$$y_{12 \dots n} = y'_{12 \dots n} \cdot \exp -\frac{1}{2}\phi(x_1 x_2 \dots x_n), \quad (22)$$

the form of the function  $\phi$  being determined by the fact that the distribution of every array must be normal, and that the mean of the array of any one variable associated with given types of the others must be the linear function of those types given by the general regression equation of the form (1). We must have, accordingly

$$\begin{aligned} \phi &= \frac{x_1^2}{\sigma_{1 \cdot 23 \dots n}^2} + \frac{x_2^2}{\sigma_{2 \cdot 13 \dots n}^2} + \dots + \frac{x_n^2}{\sigma_{n \cdot 12 \dots (n-1)}^2} \\ &- 2r_{12 \cdot 3 \dots n} \frac{x_1 x_2}{\sigma_{1 \cdot 23 \dots n} \sigma_{2 \cdot 13 \dots n}} - \dots - 2r_{(n-1) \cdot n \cdot 12 \dots (n-2)} \frac{x_{n-1} x_n}{\sigma_{(n-1) \cdot 1 \dots (n-2)} \sigma_{n \cdot 1 \dots (n-1)}}. \end{aligned} \quad (23)$$

But this expression may be thrown into several different forms. Thus, replacing the correlated variables,  $x_1, x_2, \dots, x_n$ , by the independent variables,  $x_1, x_{2 \cdot 1}, x_{3 \cdot 12}, \dots, x_{n \cdot 1 \dots (n-1)}$ , we have the very useful form

$$\phi = \frac{x_1^2}{\sigma_1^2} + \frac{x_{2 \cdot 1}^2}{\sigma_{2 \cdot 1}^2} + \frac{x_{3 \cdot 12}^2}{\sigma_{3 \cdot 12}^2} + \dots + \frac{x_{n \cdot 1 \dots (n-1)}^2}{\sigma_{n \cdot 1 \dots (n-1)}^2}. \quad (24)$$

This expression may be shown to be identical with (23) by expanding in terms of deviations of order zero, and reducing the coefficients of the square terms by means of the equation

$$\frac{r_{n \cdot k}^2}{\sigma_{n \cdot k}^2} + \frac{1}{\sigma_n^2} = \frac{1}{\sigma_{n \cdot k}^2},$$

and those of the product terms by an equation derived at once from (19),

$$\frac{r_{13 \cdot 2k} r_{23 \cdot 1k}}{\sigma_{1 \cdot 3k} \sigma_{2 \cdot 3k}} - \frac{r_{12}}{\sigma_{1 \cdot k} \sigma_{2 \cdot k}} = -\frac{r_{12 \cdot 3k}}{\sigma_{1 \cdot 3k} \sigma_{2 \cdot 3k}}.$$

21. Several important results follow at once from the form of the expression (24) for the exponent  $\phi$ . Since the variables are independent, the central value of the normal function,  $y'_{12\dots n}$ , must be given by the product of the well-known expressions for the single variables, *i.e.*, we must have

$$y'_{12\dots n} = N/(2\pi)^{n/2} \sigma_1 \sigma_{2.1} \sigma_{3.12} \dots \sigma_{n.1\dots(n-1)}. \quad (25)$$

22. Again, if we integrate the normal function in the form (24) with respect to  $x_1$ , treating the remaining variables,  $x_{2.1}$ ,  $x_{3.21}$ , etc., as constants of integration,  $\sigma_1$  is eliminated from  $y'$  and  $x_1$  from  $\phi$ , and all the remaining variables in the exponent contain the secondary suffix 1. If  $x_{2.1}$ ,  $x_{3.21}$ ,  $\dots x_{n.1\dots(n-1)}$  are then replaced by  $x_{2.1}$ ,  $x_{3.1}$ ,  $\dots x_{n.1}$ ,  $\phi$  may be written in the form (23) for these variables. Similarly, if we integrate again with respect to  $x_{2.1}$ ,  $\sigma_{2.1}$  is removed from  $y'$  and  $x_{2.1}$  from  $\phi$ , and all the remaining variables in the exponent contain both secondary suffixes 1 and 2. If  $x_{3.21}$ ,  $x_{4.321}$ ,  $\dots x_{n.321}$  are then replaced by  $x_{3.21}$ ,  $x_{4.21}$ ,  $\dots x_{n.21}$ ,  $\phi$  may be written in the form (23) for these variables. Clearly the process may be continued on the same lines. The correlation between all sets of deviations, of any one order, with the same secondary suffixes, is therefore normal correlation.

23. It follows that we may generalise at once the known formulæ for the probable errors of the constants of a normal distribution. Omitting the factor 0·674489...we have, standard error of a

Standard deviation $\sigma_1$ .....	$\sigma_1/\sqrt{2N}$ .
Correlation coefficient $r_{12}$ .....	$1 - r_{12}^2/\sqrt{N}$ .
Regression coefficient $b_{12}$ .....	$\sigma_{1.2}/\sigma_2\sqrt{N}$ .

The first is a well-known result; the last two are cited from the valuable memoir by Professor Karl Pearson and Mr. L. N. G. Filon.\* But since  $\sigma_{1.k}$  is the standard deviation of the normally distributed variable  $x_{1.k}$ ,  $r_{12.k}$  the correlation between the normally distributed variables  $x_{1.k}$  and  $x_{2.k}$ , and  $b_{12.k}$  the regression of  $x_{1.k}$  on  $x_{2.k}$ , we must have, quite generally,  $k$  denoting as before either a single subscript or an aggregate, standard error of a

Standard deviation $\sigma_{1.k}$ .....	$\sigma_{1.k}/\sqrt{2N}$ .
Correlation coefficient $r_{12.k}$ .....	$1 - r_{12.k}^2/\sqrt{N}$ .
Regression coefficient $b_{12.k}$ .....	$\sigma_{12.k}/\sigma_{2.k}\sqrt{N}$ .

(26)

The last result may be readily verified against the formula arrived at by Professor Pearson and Mr. Filon, for the case of three variables, after pages of the most laborious work.† The first may be checked for the case of two

\* 'Phil. Trans.,' A (1898), vol. 191, p. 229.

† *Loc. cit.*, equation xxxviii, p. 260.

variables, remembering the result of the same writers,\* that the correlation between errors in  $\sigma_1$  and in  $r_{12}$  is  $r_{12}/\sqrt{2}$ ; for we have

$$\sigma_{1.2} = \sigma_1 (1 - r_{12}^2)^{\frac{1}{2}};$$

$$\frac{d\sigma_{1.2}}{\sigma_{1.2}} = \frac{d\sigma_1}{\sigma_1} - \frac{r_{12} \cdot dr_{12}}{1 - r_{12}^2}.$$

Or, squaring both sides of the equation and summing, using  $\epsilon_{1.2}$  to denote the standard error of  $\sigma_{1.2}$ ,

$$\frac{\epsilon_{1.2}^2}{\sigma_{1.2}^2} = \frac{1}{2N} + \frac{r_{12}^2}{N} - \frac{2r_{12}}{1 - r_{12}^2} \frac{r_{12}}{\sqrt{2}} \frac{1}{\sqrt{2N}} \frac{1 - r_{12}^2}{\sqrt{N}} = \frac{1}{2N}.$$

23. The question of errors of sampling in the case of the coefficient of  $n$ -fold correlation,  $R$ , is not so simple, owing to the fact that the sign of the coefficient is essentially positive and, consequently, it is subject to biassed error. If, for instance, a series of variables are strictly independent, but values are found for  $r_{12}, r_{13.2}, r_{14.32}$ , etc., equal to  $\delta_2, \delta_3, \delta_4, \dots$  then

$$1 - R^2_{1(23\dots n)} = (1 - \delta_2^2)(1 - \delta_3^2) \dots (1 - \delta_n^2).$$

If the  $\delta$ 's are sufficiently small to enable us to neglect terms of the fourth order as compared with those of the second order, then we may write to the first approximation,

$$R^2_{1(23\dots n)} = \delta_2^2 + \delta_3^2 + \dots + \delta_n^2.$$

Or, summing for a number of samplings and substituting  $1/N$  for  $\Sigma(\delta^2)$  in each case, the root-mean-square value of  $R$  when the variables are strictly independent is

$$R_0 = (n-1)^{\frac{1}{2}} / N^{\frac{1}{2}}, \quad (27)$$

$n$  being the number of variables and  $N$  the number of observations.  $R$  cannot be held with certainty to be of definite significance if not markedly greater than this, and if the number of observations be small compared with the number of variables, the critical value is rather unpleasantly large. Thus in the case of a recent investigation by Mr. R. H. Hooker into the relation between the weather and the crops,  $n = 3$ ,  $N = 21$ , consequently  $R_0 = \sqrt{\frac{2}{21}} = 0.31$  (the value cited by him on my authority).† Clearly, if the number of observations be small, it is not worth while dealing with a large number of variables.

\* *Loc. cit.*, equation xviii, p. 242.

† 'Roy. Stat. Soc. Journ.', vol. 70 (1907), p. 7.